

1 Preface

This document explains how to extend the model described in Qian [2020] to cover the case where the average gene expressions (grouped by the cell class) in the single cell data are random. We also assume that the cell class prior probability π_k is also unknown and random.

2 The variational distribution

The original joint model as described in Qian is:

$$\begin{aligned} \log p(\mathbf{x}, \boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\zeta}, \mathbf{z}) = & - \sum_{g,c,k} \zeta_{c,k} \int \mu_{g,k} e^{-D_c(t)} \gamma_{g,c} \eta_g dt + \\ & + \sum_{s,c} z_{s,c} \left[-D_c(x_s) + \log(\gamma_{g_s,c}) + \log(\eta_{g_s}) + \sum_k \zeta_{c,k} \log(\mu_{g_s,k}) \right] + \\ & + \sum_{g,c} \log p(\gamma_{g,c}) + \sum_g \log p(\eta_g) + \\ & + \sum_{c,k} \zeta_{c,k} \log \pi_k \end{aligned} \quad (1)$$

2.1 The prior on $\mu_{g,k}$

To introduce a random $\mu_{g,k}$ we assume that

$$\mu_{g,k} \sim \Gamma(mM_{g,k}, m)$$

As a reminder we note that

$$\bar{\mu}_{g,k} = E[\mu_{g,k}] = M_{g,k}$$

and

$$\overline{\log \mu_{g,k}} = E[\log \mu_{g,k}] = \psi(mM_{g,k}) - \log(m)$$

2.2 The prior on π_k

Let $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_K)$ be the vector of cell class probabilities. We assume a Dirichlet prior distribution on $\boldsymbol{\pi}$:

$$(\pi_1, \pi_2, \dots, \pi_K \mid \delta_1, \delta_2, \dots, \delta_K) \sim D(\delta_1, \delta_2, \dots, \delta_K)$$

and

$$p(\pi_1, \pi_2, \dots, \pi_K) \propto \prod_{i=1}^K \pi_i^{\delta_i-1}$$

We remind that

$$\bar{\pi}_k = E[\pi_k] = \frac{\delta_k}{\sum_{i=1}^K \delta_i}$$

and

$$\overline{\log \pi_k} = E[\log \pi_k] = \psi(\delta_k) - \psi\left(\sum_{i=1}^K \delta_i\right)$$

2.3 Main

The (new) joint model is:

$$\begin{aligned} \log p(\mathbf{x}, \boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\zeta}, \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\pi}) = & - \sum_{g,c,k} \zeta_{c,k} \mu_{g,k} A_c \gamma_{g,c} \eta_g + \\ & + \sum_{s,c} z_{s,c} \left[-D_c(x_s) + \log(\gamma_{g_s,c}) + \log(\eta_{g_s}) + \sum_k \zeta_{c,k} \log(\mu_{g_s,k}) \right] + \\ & + \sum_{g,c} \log p(\gamma_{g,c}) + \sum_g \log p(\eta_g) + \sum_{g,k} \log p(\mu_{g,k}) \\ & + \sum_{c,k} \zeta_{c,k} \log \pi_k + \sum_k (\delta_k - 1) \log \pi_k \end{aligned} \quad (2)$$

where $A_c = \int e^{-D_c(t)} dt$.

2.3.1 The variational distribution of $\mu_{g,k}$

The variational distribution is:

$$\begin{aligned} \log q(\boldsymbol{\mu}) = & E_{-\boldsymbol{\mu}} \left[\log(\mathbf{x}, \boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\zeta}, \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\pi}) \right] = \\ & - \sum_{g,c,k} \bar{\zeta}_{c,k} \mu_{g,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + \sum_{s_g,c} \bar{z}_{s_g,c} \left[\sum_k \bar{\zeta}_{c,k} \log(\mu_{g_s,k}) \right] + \sum_{g,k} \log p(\mu_{g,k}) \end{aligned} \quad (3)$$

It is worth reminding that

$$z_{s_g,c} = \begin{cases} 1 & \text{if spot } s_g \text{ from gene } g \text{ belongs to cell } c \\ 0 & \text{otherwise} \end{cases}$$

thus adding together all spots from the same gene yields the total gene counts of gene g for any given cell c . We denote this total gene count by $N_{g,c}$ and:

$$N_{c,g} = \sum_{s:s_g=g} z_{s,c}$$

The quantity $\bar{z}_{s_g,c}$ is the expectation of the indicator function defined above and is also equivalent to the probability that spot s_g belongs to cell c . We denote the sum of all those probabilities by $\bar{N}_{c,g}$ and it expresses the expected counts of gene g for cell c . Hence:

$$\bar{N}_{c,g} = \sum_{s:s_g=g} \bar{z}_{s,c}$$

In equation 3 the summation over spots can now be expressed in terms of a summation over the genes:

$$\begin{aligned} \log q(\mu) &= E_{-\mu} [\log(\mathbf{x}, \boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\zeta}, \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\pi})] = \\ &= - \sum_{g,c,k} \bar{\zeta}_{c,k} \mu_{g,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + \sum_{g,c} \bar{N}_{c,g} \left[\sum_k \bar{\zeta}_{c,k} \log(\mu_{g,k}) \right] + \sum_{g,k} \log p(\mu_{g,k}) = \\ &= - \sum_{g,c,k} \bar{\zeta}_{c,k} \mu_{g,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + \sum_{g,c,k} \bar{N}_{c,g} \bar{\zeta}_{c,k} \log(\mu_{g,k}) + \sum_{g,k} \log p(\mu_{g,k}) = \\ &= \sum_{g,k} \underbrace{\left[- \sum_c \bar{\zeta}_{c,k} \mu_{g,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + \sum_c \bar{\zeta}_{c,k} \bar{N}_{c,g} \log(\mu_{g,k}) + \log p(\mu_{g,k}) \right]}_{\log q(\mu_{g,k})} \end{aligned} \quad (4)$$

If

$$\mu_{g,k} \sim \Gamma(mM_{g,k}, m) \quad (5)$$

then

$$\log p(\mu_{g,k}) \propto (mM_{g,k} - 1) \log(\mu_{g,k}) - m\mu_{g,k}$$

and from equation 4 we have:

$$\begin{aligned} \log q(\mu_{g,k}) &= - \sum_c \bar{\zeta}_{c,k} \mu_{g,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + \sum_c \bar{\zeta}_{c,k} \bar{N}_{c,g} \log(\mu_{g,k}) + \\ &\quad + (mM_{g,k} - 1) \log(\mu_{g,k}) - m\mu_{g,k} = \\ &= \left[\sum_c \bar{\zeta}_{c,k} \bar{N}_{c,g} + mM_{g,k} - 1 \right] \log(\mu_{g,k}) - \left[\sum_c \bar{\zeta}_{c,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + m \right] \mu_{g,k} \end{aligned} \quad (6)$$

Hence:

$$\mu_{g,k} \sim \Gamma\left(\sum_c \bar{\zeta}_{c,k} \bar{N}_{c,g} + m M_{g,k}, \quad \sum_c \bar{\zeta}_{c,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + m\right) \quad (7)$$

If the view expressed by the assignment probabilities $\bar{\zeta}_{c,k}$ is very weak and can be neglected then the (posterior) mean of $\mu_{g,k}$ will tend to $M_{g,k}$:

$$\bar{\mu}_{g,k} \rightarrow M_{g,k}$$

which is the a-priori and known specification for the mean counts of gene g within class k .

If on the other hand, the assignment probabilities reflect a very strong belief then:

$$\bar{\mu}_{g,k} \rightarrow \frac{\sum_c \bar{\zeta}_{c,k} \bar{N}_{c,g}}{\sum_c \bar{\zeta}_{c,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g}$$

The numerator involves the product of the expected gene counts of any given gene g in some cell c by the probability that the cell's class is k . When this product is summed over c it gives us what can be thought of as the model-fitted mean gene counts of g within class k (scaled by the denominator)

In general, the mean of the distribution in 7 is

$$\bar{\mu}_{g,k} = \frac{\sum_c \bar{\zeta}_{c,k} \bar{N}_{c,g} + m M_{g,k}}{\sum_c \bar{\zeta}_{c,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + m} \quad (8)$$

and the mean of its log-transformation:

$$\begin{aligned} \overline{\log \mu_{g,k}} &= E[\log \mu_{g,k}] \\ &= \psi\left(\sum_c \bar{\zeta}_{c,k} \bar{N}_{c,g} + m M_{g,k}\right) - \log\left(\sum_c \bar{\zeta}_{c,k} A_c \bar{\gamma}_{g,c} \bar{\eta}_g + m\right) \end{aligned} \quad (9)$$

For our studies we assume that the hyperparameter m of the prior specified in equation 5 is equal to one:

$$m = 1$$

In practice that means that the prior distribution of $\mu_{g,k}$ will yield the same mean as the variance. Also if $m < 1$ then the prior variance widens otherwise it will contract.

2.3.2 The variational distribution of π_k

$$\begin{aligned}
\log q(\boldsymbol{\pi}) &= \mathbb{E}_{-\pi} [\log(\mathbf{x}, \boldsymbol{\gamma}, \boldsymbol{\eta}, \boldsymbol{\zeta}, \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\pi})] = \\
&= \sum_{c,k} \bar{\zeta}_{c,k} \log \pi_k + \sum_k (\delta_k - 1) \log \pi_k = \\
&= \sum_k \sum_c \bar{\zeta}_{c,k} \log \pi_k + \sum_k (\delta_k - 1) \log \pi_k = \\
&= \sum_k \left[\sum_c \bar{\zeta}_{c,k} \log \pi_k + (\delta_k - 1) \log \pi_k \right] = \\
&= \sum_k \left[\left(\sum_c \bar{\zeta}_{c,k} + \delta_k - 1 \right) \log \pi_k \right] \tag{10}
\end{aligned}$$

Thus the posterior of $\boldsymbol{\pi}$ follows a Dirichlet distribution:

$$\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_K) \sim D\left(\sum_c \bar{\zeta}_{c,k=\text{class } 1} + \delta_1, \dots, \sum_c \bar{\zeta}_{c,k=\text{class } K} + \delta_K\right) \tag{11}$$

The posterior mean of the marginal π_k is

$$\bar{\pi}_k = \mathbb{E}[\pi_k] = \frac{\sum_c \bar{\zeta}_{c,k} + \delta_k}{\sum_k [\sum_c \bar{\zeta}_{c,k} + \delta_k]} \tag{12}$$

The expectation $\bar{\zeta}_{c,k}$ is the probability that cell c has class k . Hence the sum

$$\sum_c \bar{\zeta}_{c,k}$$

represents the total number of cells with class k . Thus, from equation 12, the more often a cell type appears the greater the weight it will have in the joint probability model described in equation 2, something that sounds reasonable. We note that the posterior mean of the log marginal is:

$$\overline{\log \pi_k} = \mathbb{E}[\log \pi_k] = \psi\left(\sum_c \bar{\zeta}_{c,k} + \delta_k\right) - \psi\left(\sum_k \left[\sum_c \bar{\zeta}_{c,k} + \delta_k\right]\right)$$